

- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15

^a Institute of Neuroscience and Psychology, University of Glasgow, Scotland G12
8QB, United Kingdom.

Corresponding author

*Philippe G. Schyns

Tel.: +44 (0) 141 330 4937

E-mail: philippe.schyns@glasgow.ac.uk

Corresponding author

*Philippe G. Schyns

Tel.: +44 (0) 141 330 4937

E-mail: philippe.schyns@glasgow.ac.uk

16 Current cognitive theories are cast in terms of information processing mechanisms
17 that use mental representations [1-4]. For example, people use their mental
18 representations to identify familiar faces under various conditions of pose,
19 illumination and ageing, or to draw resemblance between family members. Yet, the
20 actual information contents of these representations are rarely characterized, which
21 hinders knowledge of the mechanisms that use them. Here, we modelled the 3D
22 representational contents of 4 faces that were familiar to 14 participants as work
23 colleagues. The representational contents were created by reverse correlating
24 identity information generated on each trial with judgments of the face's similarity to
25 the individual participant's memory of this face. In a second study, testing new
26 participants, we demonstrated the validity of the modelled contents using everyday
27 face tasks that generalize identity judgments to new viewpoints, age and sex. Our
28 work highlights that such models of mental representations are critical to
29 understanding generalization behavior and its underlying information processing
30 mechanisms.

The cognitive mechanism of recognition is guided by mental representations that are stored in memory [1-4]. Personal familiarity with faces (e.g. as family members, friends or work colleagues) provides a compelling everyday illustration because the information contents representing familiar faces in memory must be sufficiently detailed to enable accurate recognition (i.e. identifying 'Mary' amongst other people) and sufficiently versatile to enable recognition across diverse common tasks—e.g. identifying Mary in different poses, at different ages or identifying her brother based on family resemblance [5-7]. And yet, it remains a fundamental challenge to reverse engineer the participant's memory to model and thereby understand the detailed contents of their representations of familiar faces. This challenge is a cornerstone to understand the brain mechanisms of face identification, because they process the contents to predict the appearance of the familiar face of 'Mary' in the visual array and to selectively extract its identity information to generalize behavior across common tasks.

We studied how our own work colleagues recognize the faces of other colleagues from memory. The work environment provides a naturally occurring and common medium of social interactions for all participants, who had at a minimum six months of exposure with the people whose faces the study tested. To model the 3D face identity information stored in their memory, we developed a methodology based on reverse correlation (see Figure 1A, and *Methods, Reverse Correlation Experiment*) and a new Generative Model of 3D Face Identity (i.e. GMF, see Figure 1B, and *Methods, Generative Model of Face Identity*), separately for 3D shape and 2D texture information (see Supplementary Figure 1A for 3D face parameters).

On each experimental trial, our GMF synthesized a set of 6 new 3D faces (see Random Faces in Figure 1A), each with a unique and randomly generated identity. Critically, each face shared other categorical face information (i.e. sex, age and ethnicity) with one of the four faces that were personally familiar to each one of our 14 participants as work colleagues—e.g. the familiar target face of 'Mary'. To achieve this, we used a General Linear Model (GLM) to decompose the familiar target face into a categorical component (e.g., for 'Mary' the average of all white females faces of 30 years of age) plus a residual component that defines the specific identity of the familiar face (see *Identity Modelling* in Figure 1B). We then generated new random identities by keeping the categorical component of the target constant (e.g., white female, 30 years of age) and adding a random component of identity (see *Identity Generation* in Figure 1B, and *Methods, Reverse Correlation Experiment, Random Face Identities* for details). Participants saw these randomly generated faces in full frontal view and selected the one that most resembled the familiar target (e.g., 'Mary') and rated its similarity to the target on a 6-point Likert scale, ranging from not at all ('1') to highly similar ('6'). To resolve the task, participants must compare the randomly generated faces presented on each trial with their mental representation of the familiar target in full frontal view. Therefore, each face selected comprises a match to the participant's mental representation of the target, which is estimated by the similarity rating of that face.

After many such trials, we used reverse correlation [8] to estimate the information content of the mental representation of each target familiar face ($N = 4$, see Supplementary Figure 1B) in each participant ($N = 14$, see *Methods, Reverse Correlation Experiment*). Specifically, we build a statistical relationship between the information content of the faces that the participant selected on each trial with their corresponding similarity ratings. In a second stage, we tested with a new group of participants ($N = 12$, i.e. the validators, see *Methods, Generalization Experiments*) whether these modelled mental representations were sufficiently detailed to enable identification of each target familiar face and sufficiently versatile to enable resemblance judgments across diverse everyday tasks—i.e. generalization across new viewpoints, age and siblings.

To reconstruct the information contents of mental representations, we used linear regression to compute the single-trial relationship between <similarity ratings, random face identity components> for each target familiar face and participant. Specifically, we computed separate regressions between the similarity ratings and each 3D shape vertex and each RGB texture pixel that comprise the face identity components. We then used the resulting Beta coefficients to model the 3D shape and texture identity components that characterize the participant's mental representation of each familiar face in the GMF (see Supplementary Figure 2 and *Methods, Analyses, Linear Regression Model and Reconstructing Mental Representations*).

With this approach, we can formally characterize and then compare the participant's mental representation of a familiar face with the ground truth face—i.e. the objective identity component of the scanned familiar face, see Supplementary Figure 1B. We focus only on 3D shape because there were very few and non-systematic relationships for texture (see Supplementary Figure 3). To illustrate, grey faces on the x-axis of Figure 2A show the ground truth identity component of 'Mary' in the GMF for Inward and Outward 3D shape deviations in relation to the categorical average (i.e., of all white females of 30 years of age, like 'Mary'). For example, Mary's nose is objectively thinner than the average of white females of her age, and so these vertices deviate inward (darker grey tones indicate increasing deviations). Likewise, her more pouty mouth is shown as an outward 3D shape deviation. The y-axis of Figure 2A uses the same format to show the mental representation of Mary in one typical participant, where colors indicate increasing deviations. These contents reveal faithful representations of, for example, a thinner nose and a pouty mouth (see *Methods, Analyses, Vertex Contribution to Mental Representations*). A scatter plot visualizes the vertex by vertex fit between the mental representation (y-axis) and the ground truth 3D face (x-axis). The white diagonal line provides a veridical reference, where the identity component in the mental representation is identical to the ground truth face, for every single 3D vertex. This is because the mental representation and ground truth faces are both registered in the same space of 3D vertices [9].

115 Our analyses reveal the specific vertices near the veridical line that faithfully
116 represent 'Mary' in the mind of this participant as colored dots reported on the scatter
117 and located on the y-axis faces in Figure 2A ($p < 0.05$, two-sided, using a random
118 permutation test to generate a chance level distribution, as reported in *Methods,*
119 *Analyses, Vertex Contribution to Mental Representation*). In contrast, white vertices
120 away from the veridical line did not faithfully represent the identity. We repeated the
121 analysis of represented contents for each participant ($N = 14$) and familiar face ($N =$
122 4). Figure 2B reports the collated group results, using the format of Figure 2A, where
123 colors now indicate N , i.e. the number of participants who faithfully represented that
124 identity in their mind with this particular 3D shape vertex. Figure 2B demonstrates
125 that mental representations comprised similar information contents across the 14
126 individual participants. Most (10/14) faithfully represented 'Mary's' thin nose, 'John's'
127 receding eyes and wider upper face (13/14), 'Peter's' prominent eyebrow and jawline
128 (13/14), 'Stephany's' protruding mouth (13/14).

129 Such convergence of represented contents across participants suggests that
130 the face representations could be multivariate (i.e. comprising contiguous surface
131 patches rather than isolated vertices). As a final step, we extracted the main
132 multivariate components of represented surface patches. To this end, we applied
133 across observers ($N = 14$) and familiar faces ($N = 4$) the Non-negative Matrix
134 Factorization (NNMF, [10]) to the faithfully represented 3D vertices (see *Methods,*
135 *Analyses, Components of Memory Representation*). Figure 3A shows the multivariate
136 components that faithfully represent four target identities and Figure 3B shows their
137 combinations for the diagnostic components of each target identity (e.g. for 'Mary,'
138 the red background heatmap; for 'Stephany,' the green one and so forth). Importantly,
139 these diagnostic components of familiar face identity have complementary
140 nondiagnostic components (i.e. the grey background heatmaps in Figure 3B), which
141 capture variable face surfaces that do not comprise the participants' mental
142 representations.

143 Here, we develop the critical demonstration that the information contents of
144 the mental representations we modelled are valid. That is, the contents enable
145 accurate identification of each target face and they also enable resemble tasks that
146 preserve their identity. We asked a new group of participants (called 'validators') to
147 resolve a variety of resemblance tasks that are akin to everyday tasks of face
148 recognition. Success on these tasks would demonstrate that the diagnostic
149 components derived from the previous experiment comprise identity information that
150 can be used in a different generalization tasks. Therefore, although the components
151 are extracted under one viewpoint (full-face), one age (for each identity) and one sex
152 (that of the identity), here we tested the generalization of identification performance
153 to new viewpoints, ages and sex.

154 For this demonstration, we synthesized new diagnostic (vs. nondiagnostic)
155 faces that were parametrically controlled for the relative strength of the diagnostic
156 multivariate components of identity vs. their nondiagnostic complement (see Figure

157 4A and *Methods, Generalization Experiments, Stimuli*). It is important to emphasize
158 that both diagnostic and nondiagnostic faces are equally faithful representations of
159 the original ground truth. That is, their shape features are equidistant from the shared
160 categorical average. However, whereas the diagnostic components deviate from the
161 average with multivariate information extracted from the participants' mental
162 representations, the nondiagnostic components do not. We hypothesized that,
163 though equidistant from the categorical average, only the diagnostic components will
164 impact performance on the resemblance tasks. For all synthesized faces, we
165 changed their viewpoint (rotation of -30 deg, 0 deg and +30 deg in depth), age (to 80
166 years old), and sex (to opposite) using the generative model--see Supplementary
167 Figure 5 to 8 for each familiar target.

168 In three independent resemblance tasks – changes of viewpoint, age and sex
169 – we tested the identification performance of 12 validators on the diagnostic and
170 nondiagnostic faces using a 5 Alternative Force Choice task (i.e. responding one of
171 four familiar identities plus a 'don't know' response, see *Methods, Generalization*
172 *Experiments, Procedure*). In each task, for each identity we found a significantly
173 higher identification performance for diagnostic faces (see Figure 4B, red curves)
174 than for nondiagnostic faces (black curves)—i.e. a fixed effect of Face Type in a
175 mixed effects linear model. For 'Mary', $F(1, 12.76) = 315.49, p < 0.001$, estimated
176 slope = 0.297, 95% Confidence Intervals = [0.264, 0.33]; for 'Stephany', $F(1, 20.62)$
177 = 25.068, $p < 0.001$, estimated slope = 0.058, 95% Confidence Intervals = [0.035,
178 0.081]; for 'John', $F(1, 12) = 21.369, p < 0.001$, estimated slope = 0.143, 95%
179 Confidence Intervals = [0.083, 0.204]; for 'Peter', $F(1, 12.01) = 5.76, p = 0.034$,
180 estimated slope = 0.095, 95% Confidence Intervals = [0.017, 0.173] (see *Methods,*
181 *Generalization Experiments, Analyses* for the detailed specification and
182 Supplementary Table 3 to 6 for the full statistical analysis of the models). Thus, the
183 diagnostic contents of the mental representations we modelled do indeed contain the
184 information that can resolve identity and resemblance tasks.

185 Mental representations stored in memory are critical to guide the information
186 processing mechanisms of cognition. Here, with a methodology based on reverse
187 correlation and a new 3D face information generator (i.e. our 3D GMF), we modelled
188 the information contents of mental representations of 4 familiar faces in 14 individual
189 participants. We showed that the contents converged across participants on a set of
190 multivariate features (i.e. local and global surface patches) that faithfully represent
191 3D information that is objectively diagnostic of each familiar face. Critically, we
192 showed that validators could identify new faces generated with these diagnostic
193 representations across three resemblance tasks—i.e. changes of pose, age and
194 sex—but performed much worse with equally faithful, but nondiagnostic features.
195 Together, our results demonstrate that the modelled representational contents were
196 both sufficiently precise to enable face identification within task and versatile enough
197 to generalize usage of the identity contents to other resemblance tasks.

198 At this stage, it worth stepping away from the results and emphasize that it is
199 remarkable that the reverse correlation methodology works at all, let alone produce
200 robust generalization across resemblance tasks. In the experiment, we asked
201 observers to rate the resemblance between a remembered familiar face, and
202 randomly generated faces, that by construction are very unlike the target face (never
203 identical, and almost never very similar). And yet, our results show that the
204 representational contents we modelled following such a task were in fact part of the
205 contents that objectively (i.e. faithfully) support identity recognition. This raises a
206 number of important points that we now discuss.

207 There has been a recent surge of interest in modelling face representations
208 from human memory [11-13]. These studies used 2D face images and applied
209 dimensionality reduction (e.g. PCA [14] and multidimensional scaling) to formalize an
210 image-based face space, where each dimension is a 2D eigenface or classification
211 image – i.e. pixel-wised RGB (or L*A*B) values. To understand the contribution of
212 each 2D face space dimension to memory representations (including their neural
213 coding), researchers modelled the relationship between projected weights of the
214 original 2D face images on each dimension and participants' corresponding
215 behavioral [13] (and brain [11, 12]) responses.

216 These studies contributed important developments in face identification
217 research because they addressed the face identity contents that the brain uses to
218 guide face identification mechanisms. Our aim was to model the face identity
219 contents in the generative 3D space of faces (not the 2D space of their image
220 projections) and to use these models to generate identification information in
221 resemblance tasks that test the generalizability of identity information. It is important
222 to clarify that we modelled identity information in a face space that belongs to the
223 broad class of 3D morphable, Active Appearance Models of facial synthesis (AAMs,
224 [15, 16]). These models contain full 3D surface and 2D texture information about
225 faces and so with their better control superseded the former generation of 2D image-
226 based face spaces ([14, 17] [18]). To synthesize faces, we used our GMF to
227 decompose each face identity as a linear combination of components of 3D shape
228 and 2D texture added to a local average (that summarizes the categorical factor of
229 age, gender, ethnicity and their interactions, cf. Figure 1B). To model the mental
230 representations of faces, we estimated the identity components of shape and texture
231 from the memory of each observer. These components had generative capacity and
232 we used them to precisely control the magnitude of identity information in new faces
233 synthesized to demonstrate generalization across pose, age and sex. Thus, we used
234 the same AAM framework for stimulus synthesis, mental representation estimation
235 and generation of generalizable identities.

236 There is a well-known problem with using AAMs to model the psychology of
237 face recognition. Perceptual expertise and familiarity are thought to involve
238 representations of faces that enable the greater generalization performance that is
239 widely reported [19-22]. However, AAMs typically adopt a brute force approach to

identity representation: a veridical (i.e. totally faithful) deviation of each physical shape vertex and texture pixel from an average. Thus, as AAMs overfit identity information, they appear as a priori weak candidate models to represent perceptual expertise with faces [18]. Our approach of studying the contents of mental representations suggests a solution to this conundrum. We showed that each observer faithfully represented only a proportion of the objective identity information that defines a familiar face identity. Our key theoretical contribution to face space is to formalize the subjective 3D diagnostic information as a reduced set of multivariate face features that can be construed as dimensions of the observer's face space. Observers develop these dimensions when they interact with the objective information that represents a new face identity in the real world. We modelled the objective information that is available to the observer for developing their face space dimensions via learning as the veridical shape and texture information of the AAM [18, 23, 24]. Key to demonstrating the psychological relevance of our psychological 3D face space dimensions is that they should comprise identity information sufficiently detailed to enable accurate face identification and sufficiently versatile to enable similarity judgments of identity in resemblance tasks. We demonstrated this potential when validators identified faces synthesized with the diagnostic dimensions in novel resemblance tasks. Thus, by introducing reduced faithful mental representations of identity information in the objective representations of AAMs we provide the means of modelling the subjective psychological dimensions of an individual's face space.

Our work could be extended to precisely track the development of the psychological dimensions of face space if we tasked observers with learning new identities (an everyday perceptual expertise task [18, 25]). Our AAMs enable a tight control of objective face information at synthesis, such as ambient factors of illumination, pose and scale, but also categorical factors of gender, sex, age and ethnicity and components of identity. Thus, we could tightly control the statistics of exposure to faces in individual observers (even orthogonalize them across observers), and model and compare the diagnostic dimensions of the psychological face space that are learned, and finally test their efficacy as we did here. And when we understand how ambient and categorical factors influence performance as a function of differential perceptual learning, we can switch to understanding familiar face identification in the wild, by progressively introducing simulations of ambient factors (e.g. identifying the face of someone walking by a street lamp at night) and observe their specific effects on performance (e.g. ambient changes in face size, shading, and cast shadows). Otherwise, all ambient and categorical factors remain naturally mixed up, and the influence of each factor to identification performance becomes near impossible to disentangle, precluding a detailed information processing understanding of face identification mechanisms.

Our results could suggest that the representation of face shape information trumps its texture. At this stage, it is important to clarify that shape and texture have different meanings in different literatures. For example, some authors in psychology

283 discuss *shape-free faces* when referring to 2D images synthesized by warping an
284 identity-specific texture to an identical ‘face shape’ (defined as a unique and standard
285 set of 2D coordinates that locate a few face features [26]). However, it is important to
286 emphasize that the warped textures are not free of 3D shape information (e.g. that
287 which can be extracted from shading [27]). In computer graphics, the generative
288 model of a face comprises a 3D shape per identity (here, specified with 4,735 3D
289 vertex coordinates), lighting sources (here, $N = 4$), and a shading model (here,
290 Phong shading [28]). The shading model interacts with shape and texture to render
291 the 3D face as a 2D image. To illustrate the effects of this rendering, Supplementary
292 Figure 9 shows how applying the same 2D textures (rows) to different 3D face
293 shapes (columns) generates 2D images with different identities. We used the better
294 control afforded by computer graphics to generate our face images and found that
295 shaded familiar face shape was more prevalent in the face memory of individual
296 participants than face texture.

297 A general question with reverse correlation tasks is whether the resulting
298 models represent a particular visual category (here, the visual identity of a face) or
299 the task from which the model was reconstructed [24, 29-31]. We contributed to this
300 debate by showing that the identity information reconstructed in one task had efficacy
301 in other tasks that involved identity. Importantly, the tasks were designed to test two
302 classes of factors: ambient and categorical. For example, we showed that the identity
303 component extracted in one ambient viewpoint (full face, 0 deg) could be used to
304 generalize identification of the same face under two new ambient viewpoints (-30 and
305 +30 deg of rotation in depth). We also showed that the identity component extracted
306 for identities (all < 40 years of age) generalized to older age (80 years). Furthermore,
307 we also showed that though extracted from a given sex, the identity component
308 would generalize to another sex, a kinship task. Hence, we found no dramatic
309 differences due to the effect of task of extraction of the identity component. Rather,
310 the extracted representational basis is useful for all tasks tested, whether using
311 ambient or categorical factors of face variance. This therefore suggests that we have
312 tapped into some essential information about familiar face representation. However,
313 we acknowledge that the generalizations we observe might still be a function of an
314 interaction between the nature of memory and the similarity task from which we
315 estimated the identity component. The component could have differed had the task
316 been more visual than memory based (e.g. identification of the same face under
317 different orientations, or a visual matching task) and we might not have derived an
318 identity component that enabled such effective generalization. In any case, the
319 memorized identity components that enable task generalization reflect an interaction
320 between memory and the input information available to represent this identity [24, 32].
321 Observers can compare this memory representation for that identity with a
322 representation of the visual input for successful identification.

323 Our models of mental representation should be construed as the abstract
324 information goals (i.e. the contents) that the visual system predicts when identifying
325 familiar faces. We call them ‘abstract information goals’ because they reflect the

326 invariant visual representations that enable the resemblance response and must be
327 broken down into global and local constituents according to the constraints of
328 representation and implementation at each level of the visual hierarchy—or their
329 analogues in deep convolutional networks, where we can use a similar methodology
330 to understand the identity contents represented in the hidden layers [33]. In norm-
331 based coding [17, 34], face identity information is represented in reference to the
332 average of a multi-dimensional face space. Monkey single cell responses increase
333 their firing rate with increasing distance of a face to this average (as happens with e.g.
334 caricaturing, [35]). As shown by Chang et al. [36], neurons selectively respond along
335 a single axis of the face space, not to other, orthogonal axes. An interesting direction
336 of research is to determine whether our reduced diagnostic features, as defined by
337 our ‘abstract information goal’ (see also [37]), provide a superior fit to the neural data
338 than the full feature sets used in the axis model used by Chang et al. [36].

339 Though we modelled the mental representation of a face identity in an AAM, it
340 is important to state that we do *not* assume that memory really represents faces in
341 this way (i.e. as demarcations to an average, separately for 3D shape and 2D
342 texture). AAM is only a state-of-the-art, mathematical modelling framework. We fully
343 acknowledge there are many possible concrete implementations into a neural, or a
344 neurally-inspired architecture that could deliver AAM-like performance without
345 assuming an explicit AAM representation. What is clear is that whichever
346 implementation, in whichever architecture, the abstract information modelled under
347 AAM framework will have to enable the performance characteristics our resemblance
348 tasks demonstrated.

349 For example, we would hypothesize that the diagnostic identity components
350 in Figure 3B are broken down, bottom to top, into the representational language of
351 V1—i.e. as representation in multi-scale, multi-orientation Gabor-like, retinotopically
352 mapped receptive fields [38, 39]; at intermediate levels of processing, as the sort of
353 local surface patches [40, 41] that we reveal, and at the top level as the combinations
354 of surface patches that enable identification and resemblance responses. Under a
355 framework of top-down prediction [42, 43], the abstract information goal of a familiar
356 face identity should trim, in a top-down manner, the fully-mapped but redundant
357 information on the retina into the task-relevant features that are transferred along the
358 occipital to ventral/dorsal visual hierarchy [37]. Tracing the construction of such a
359 reduced memory representation of face identity in the brain should enable an
360 accurate and detailed modelling of the processing mechanism along the visual
361 hierarchy (see also [12, 44-46]). What our work critically provides is an estimate of
362 the end goal of the hierarchy (i.e. the diagnostic component), which is also a
363 prediction of what is important in the input. It is in this sense that mental
364 representations guide task-specific information processing in the brain. Without
365 knowing mental representations, we do not have even have an information needle to
366 search in the fabled haystack of brain activity, let alone reconstruct the mechanisms
367 that process its contents.

368 We modelled the critical mental representations of that guide the processing
369 of visual information of familiar face identities. In several resemblance tasks that
370 require usage of face identity, we demonstrated the efficacy of the contents we
371 modelled. Our approach and results open new research avenues for the interplay
372 between visual information, categorization tasks and their implementation as
373 information processing mechanisms in the brain.

374 METHODS

375 Generative Model of 3D Face Identity (GMF).

376 We designed a generative model to objectively characterize and control 3D face
377 identity variance, using a database of 355 3D faces (acquired with a 4D face capture
378 system, see *Supplementary Methods, 3D Face Database*) that describes each face
379 by its shape (with 3D coordinates for each one of 4,735 vertices) and its texture (with
380 the RGB values of 800*600 pixels, see Supplementary Figure 1A). It is critical to
381 reiterate that the familiar faces were not part of the 3D face database.

382 To design the 3D GMF, we first applied a high-dimensional General Linear
383 Model (GLM), separately to 3D vertex coordinates and 2D pixel RGB values, to
384 model and explain away variations in face shape and texture that arise from the non-
385 identity categorical factors of sex, age, ethnicity, and their interactions. The GLM
386 therefore: 1) extracted as a non-identity face average the shape and texture face
387 information explained by non-identity categorical factors; and also 2) isolated the
388 residual information that defines the 3D shape and 2D texture identity information of
389 each face--i.e. the identity residuals.

390 To further control identity information, we applied Principal Components
391 Analysis (PCA) to the identity residuals of the 355 faces, separately for shape and
392 texture. The PCA represented shape residuals as a 355-dimensional vector in a 355-
393 dimensional space of multivariate components, and a separate PCA represented the
394 texture residuals as a 355*5 (spatial frequency bands)-dimensional matrix in a space
395 of 355*5 multivariate components. Two sets of PCA coordinates therefore
396 represented the objective shape and texture information of each identity in the
397 principal components space of identity residuals.

398 Our 3D GMF is formally expressed as follows:

$$Faces = Design\ Matrix \times Coefficient\ Matrix + weights \times PCs$$

399 Where *Faces* is the vertex (or texture) matrix of 355 faces: for vertices, it is
400 [355 x 14,205] where 14,205 = 4,735 vertices x 3 coordinates; for texture, it is [355 x
401 1,440,000] where 1,440,000 = 800 x 600 pixels x 3 RGB. *Design Matrix* defined the
402 non-identity categorical factors and their interactions (N = 9), i.e. constant, age,
403 gender, white Caucasian (WC), eastern Asian (EA), black African (BA), gender x WC,
404 gender x EA, gender x BA, for each of face (N = 355), and therefore is [355 x 9]. We
405 estimated the linear effects of each non-identity factor and their interactions using the
406 GLM which are represented in the *Coefficient Matrix* (i.e. [9 x 14,205] for shape and
407 [9 x 1,440,000] for texture). After the GLM fit, the [355 x 14,205] shape (or [355 x
408 140,000] texture) residuals are further explained using the PCA analysis, resulting
409 355 components.

410 Furthermore, Supplementary Figure 1B illustrates how the generative model
411 controlled the non-identity and identity factors using the 4 familiar faces of our

412 experiment. First, we scanned the four familiar faces of the experiment (2nd column).
413 We fitted each into our 3D GMF to derive a ground truth face (the 3rd column), with
414 minimal distortions (shown in the 1st column).

415 The model generates new 3D faces by adding the identity residuals of four
416 familiar faces to different non-identity GLM averages, to change their age, sex or
417 ethnicity separately, or jointly sex and ethnicity. The outcomes are older, sex
418 swapped, ethnicity swapped and sex and ethnicity swapped versions of the same
419 identity (the 4th to 7th column). We used these generative properties to derive the
420 stimuli of the generalization experiment.

421 **Reverse Correlation Experiment**

422 **Participants.** We recruited 14 participants (all white Caucasians, 7 females,
423 mean age = 25.86 years, SD = 2.26 years) who were personally familiar with each
424 familiar identity as work colleagues for at least 6 months. We assessed familiarity on
425 a 9-point Likert scale, from not at all familiar '1' to highly familiar '9'. Supplementary
426 Table 1 reports the familiarity ratings for each identity and participant. We chose a
427 sample size similar to those reported elsewhere [47-49]. All participants had normal
428 or corrected-to-normal vision, without a self-reported history or symptoms of
429 synaesthesia, and/or any psychological, psychiatric or neurological condition that
430 affects face processing (e.g., depression, autism spectrum disorder or
431 prosopagnosia). They gave written informed consent and received £6 per hour for
432 their participation. The University of Glasgow College of Science and Engineering
433 Ethics Committee provided ethical approval.

434 **Familiar Faces.** We scanned four faces 'Mary' and 'Stephany' (white
435 Caucasian females of 36 and 38 of age, respectively), and 'John' and 'Peter' (white
436 Caucasian males of 31 and 38 years of age, respectively) who were familiar to all
437 participants as work colleagues. As we will explain, we used these scanned faces to
438 compare the objective and mentally represented identity information in each
439 participant. Each of these four people gave informed consent for the use of their
440 faces in published papers.

441 **Random Face Identities.** We reversed the flow of computation in the 3D
442 GMF to synthesize new random identities while controlling their non-identity factors
443 (see Figure 1B *Identity Generation*, the reverse direction is indicated by the dashed
444 line). We proceeded in three steps: First, we fitted the familiar identity in the GLM to
445 isolate its non-identity averages, independently for shape and texture. Second, we
446 randomized identity information by creating random identity residuals—i.e. we
447 generated random coefficients (shape: 355; texture: 355*5) and multiplied them by
448 the principal components of residual variance (shape: 355; texture: 355*5). Finally,
449 we added the random identity residuals to the GLM averages to create a total of
450 10,800 random faces per familiar identity in the reverse correlation experiment.

451 **Procedure.** Each experimental block started with a centrally presented frontal
452 view of a randomly chosen familiar face (henceforth, the target). On each trial of the
453 block, participants viewed six simultaneously presented randomly generated
454 identities based on the target, displayed in a 2 x 3 array on a black background, with
455 faces subtending an average of 9.5° by 6.4° of visual angle. We instructed
456 participants to respond on one of 6 buttons to choose the face that most resembled
457 the target. The six faces remained on the screen until response. Another screen
458 immediately followed instructing participants to rank the similarity of their choice to
459 the target, using a 6-point Likert scale ('1' = not similar, '6' = highly similar) with
460 corresponding response buttons. Following the response, a new trial began. The
461 experiment comprised 1,800 trials per target, divided into 90 blocks of 20 trials each,
462 run over several days, for a grand total of 7,200 trials that all validators accomplished
463 in a random order. Throughout, participants sat in a dimly lit room and used a chin
464 rest to maintain a 76 cm viewing distance. We ran the experiment using the
465 Psychtoolbox for MATLAB R2012a. Data collection and following analysis were not
466 performed blind to the target faces.

467 Analyses

468 **Linear Regression Model.** For each participant and target face, each trial
469 produced two outcomes: one matrix of 4,735*3 vertex (and 800*600 RGB pixel)
470 parameters corresponding to the shape (and texture) residuals of the chosen random
471 face on this trial, and one corresponding integer that captures the similarity between
472 the random identity parameters and the target. Across the 1,800 trials per target, we
473 linearly regressed (i.e. RobustFit, Matlab 2013b) the 3D residual vertices (separately
474 for the X, Y and Z coordinates) and residual RGB pixels (separately for R, G and B
475 color channel) with the corresponding similarity rating values. These linear
476 regressions produced a linear model with coefficients Beta_1 and Beta_2 vectors for
477 each residual shape vertex coordinate and residual RGB texture pixel, for each
478 familiar face and participant. Supplementary Figure 2A illustrates the linear
479 regression model for the 3D vertices of 'Mary.' Henceforth, we focus our analyses on
480 the Beta_2 coefficients because they quantify how shape and texture identity
481 residuals deviate from the GLM categorical average to represent the identity of each
482 familiar face in the memory of each participant.

483 **Reconstructing Mental Representations.** Beta_2 coefficients can be
484 amplified to control their relative presence in a newly synthesized 3D face.
485 Supplementary Figure 2B1 illustrates such amplification for one participant's Beta_2
486 coefficients of shape and texture of 'Mary.' Following the reverse correlation
487 experiment, we brought each participant back to fine-tune their Beta_2 coefficients
488 for each familiar face, using the identical display and viewing distance parameters as
489 in the reverse correlation experiment (see Supplementary Figure 2B2 and
490 *Supplementary Methods, Fine-tuning Beta_2 Coefficients*).

491 **Vertex Contribution to Mental Representations.** Vertices, whether in the
492 ground truth face or in the participant's mental representation can deviate inward or

outward in 3D from the corresponding vertex in the common categorical average of their GLM fits (cf. Figure 1B). Thus, we can compare the respective deviations of their 3D vertices in relation to the common GLM categorical average. To evaluate this relationship, we plotted the normalized deviation of ground truth vertices from most Inward (-1) to most Outward (+1) on the X-axis of a 2D scatter plot; we also reported the normalized deviation of corresponding vertex of the mental representation on the Y-axis (as shown Figure 2A). If ground truth and mental representations were identical, their vertex-by-vertex deviations from the GLM categorical average (i.e. Euclidean distance) would be identical and would form the veridical diagonal straight white line provided as a reference in the scatter plot of Figure 2A.

Using this veridical line as a reference, for each participant and familiar face representation, we proceeded in three steps to classify each vertex as either 'faithful' or 'not faithful', and to test whether the vertices in mental representations deviated from the categorical average more than would be expected to occur by chance.

Step 1: We constructed a permutation distribution by iterating our regression analysis 1,000 times with random permutations of the choice response across the 1,800 trials. To control for multiple comparisons, we selected maximum (vs. minimum) Beta_2 coefficients across all shape vertices (and texture pixels), separately for the X, Y and Z coordinates (RGB color channels) from each iteration. We used the resulting distribution of maxima (and minima) to compute the 95% confidence interval of chance-level upper (and lower) Beta_2 value and classified each Beta_2 coefficient as significantly different from chance ($p < 0.05$, two-sided), or not. We consider the vertex (or pixel) as significant if the Beta_2 coefficient of any coordinate (or color channel) was significant. There were very few significant pixels, with almost no consistency across participants (see Supplementary Figure 3), so we excluded texture identity residuals from further analyses.

Step 2: We used the chance-fit Beta coefficients in Step 1 and the Beta_2 amplification value derived in **Reconstructing Mental Representation** to compute the equation $GLM + \beta_1 + \beta_2 * amplification\ value$ (cf. Supplementary Figure 2B). As a result, we built a distribution of 1,000 chance fit faces.

Step 3: To classify whether each significant 3D vertex in the mental representation of a participant is more similar to ground truth than we would expect by chance, we computed D_{chance} , the mean Euclidean distance between the 1,000 chance fit faces and the veridical line, and D_{memory} , the distance between the same mental representation vertex and the veridical line. If $D_{memory} < D_{chance}$, this significant vertex is 'faithful' because it is significantly closer to the veridical line than chance (and we plot it with blue to red colors in Figure 2A); if $D_{memory} > D_{chance}$, the vertex is not faithful (and we plot it in white in Figure 2A, together with the nonsignificant vertices).

533 To derive group results, we counted across participants the frequency of each
534 faithful vertex and used a Winner-Take-All scheme to determine group-level
535 consistency. For example, if 13/14 participants represented this particular vertex as
536 'faithful,' we categorized it as such at the group level and reported the number of
537 participants as a color indicating 13 participants. If there was no majority for a vertex,
538 we color-coded it as white (see Figure 2B).

539 **Components of Memory Representation.** The purpose of the following
540 analysis was to find common diagnostic components (multivariate features) that
541 emerged in the group-level memory representation of each face identity. To do so,
542 we factorized with Non-negative Matrix Factorization (NNMF) the total set of memory
543 representations across familiar identities and observers.

544 For each participant, we recoded each vertex in the identity residuals of each
545 familiar face as 'faithful' = 1, 'not faithful' or not significant = 0, resulting in a 4735-d
546 binary vector. We pooled 56 such binary vectors (across 4 targets x 14 observers =
547 56) to create a 4735 by 56 (i.e. vertex-by-model) binary matrix to which we applied
548 NNMF to derive 8 multivariate components that captured the main features that
549 faithfully represent familiar faces in memory across participants (see *Supplementary*
550 *Methods, Non-negative Matrix Factorization*). Heatmap in Figure 3A shows each
551 NNMF component.

552 To determine the loading (i.e. the contribution) of each NNMF component in
553 the group-level mental representation of each familiar face identity, we computed the
554 median loading of this component on the 14 binary vectors representing this identity
555 in the 14 observers. We applied a 0.1 loading threshold (> 73 percentile of all 8
556 components \times 4 identities median loadings) to ascribe a given component to a
557 familiar face representation. The boxplot in Figure 3A represents the loading of each
558 NNMF component at the group-level representation, with colored boxes showing at
559 least 2 above-threshold NNMF components represent each familiar identity.

560 We then constructed the diagnostic component of a familiar identity
561 representation as follows: for each vertex we extracted the maximum loading value
562 across the NNMF components representing it, and normalized the values to the
563 maximum loading across all vertices. This produced a 4735-d vector V_d that weighs
564 the respective contribution of each 3D vertex to the faithful representation of this
565 familiar identity that we call the "diagnostic component." The heat maps in the left
566 column of Figure 3B represent the diagnostic component of each familiar identity.
567 Supplementary Figure 4 shows the high accuracy of the features captured by the
568 components.

569 Crucially for our validation experiment, we were then able to define a
570 nondiagnostic component as the complement of the diagnostic component $V_n = 1 -$
571 V_d . It is important to emphasize that we adjusted the total deviation magnitude of the
572 diagnostic and nondiagnostic components from the categorical average—i.e. by
573 equating the total sum of their deviations. This ensures that diagnostic and

574 nondiagnostic components are both equidistant from the average face in the
575 objective face space. The right column of Figure 3B shows the nondiagnostic
576 component of each familiar identity representation.

577 **Generalization Experiments**

578 **Validators.** We recruited 12 further participants (7 white Caucasian and 1
579 East Asian females, 5 white Caucasian males, with mean age = 28.25 years and SD
580 = 4.11 years), using the same procedure and criteria and those presiding for the
581 selection of participants. Supplementary Table 2 reports the familiarity ratings for
582 each identity and validator. All validators had normal or corrected-to-normal vision,
583 without a self-reported history or symptoms of synaesthesia, and/or any
584 psychological, psychiatric or neurological condition that affects face processing (e.g.,
585 depression, autism spectrum disorder or prosopagnosia). They gave written informed
586 consent and received £6 per hour for their participation. The University of Glasgow
587 College of Science and Engineering Ethics Committee provided ethical approval.

588 **Stimuli.** For each familiar identity, we synthesized new 3D faces that
589 comprised graded levels of either the diagnostic or the nondiagnostic shape
590 components as explained in the section **Components of Memory Representation**
591 above. Specifically, we used the normalized diagnostic component V_d and its
592 nondiagnostic complement V_n to synthesize morphed faces with shape information of
593 each target identity as follows:

$$\text{Diagnostic Faces} = \text{Ground Truth} \times V_d \times \alpha + \text{Categorical Average} (1 - V_d \times \alpha)$$

$$\text{Nondiagnostic Faces} = \text{Ground Truth} \times V_n \times \alpha + \text{Categorical Average} (1 - V_n \times \alpha)$$

594 with amplification factor $\alpha = 0.33, 0.67, 1, 1.33, 1.67$, to control the relative
595 intensity of diagnostic and nondiagnostic shape changes. We rendered all these
596 morphed shapes with the same average texture. The first rows of Supplementary
597 Figure 5 to 8 show the morphed faces for each familiar identity. We added as filler
598 stimuli the grand average face (for both shape and texture) of the 355 database
599 faces.

600 We also changed the viewpoint, age and sex of all of these synthesized faces.
601 Specifically, we rotated them in depth by -30 deg, 0 deg and +30 deg and using the
602 3D GMF, we set the age factor to 80 years/swapped the sex factor, keeping all other
603 factors constant (cf. *Generative Model of 3D Face Identity* in Figure 1B and
604 Supplementary Figure 1B).

605 **Procedure.** The experiment comprised 3 sessions (viewpoint, age and sex)
606 that all validators accomplished in a random order, with one session per day. In the
607 Viewpoint session, validators ran 15 blocks of 41 trials (5 repetitions of 123 stimuli).
608 Each trial started with a centrally displayed fixation for 1s, followed by a face on a
609 black background for 500ms. We instructed validators to name the face as 'Mary,'
610 'Stephany,' 'John' or 'Peter,' or respond 'other' if they could not identify the face.

Validators were required to respond as accurately and as quickly as possible. A 2s fixation separated each trial. Validators could break between blocks. In the Age and Sex sessions, validators ran 5 blocks that repeated 44 trials. They were instructed to respond “Old Mary,” “Old Stephany,” “Old John,” “Old Peter” or “Other” in the age session, and “Mary’s brother,” “Stephany’s brother,” “John’s sister,” “Peter’s sister” or “Other” in the sex session. For each session, stimuli are randomized across all trials. Across the 3 sessions, we recorded participants’ identification performance in 3 viewpoints, a change of age information and a change of sex information. Data collection and following analysis were not performed blind to the conditions of the experiments.

Analyses. For each validator and generalization condition, we computed the percent correct identification of diagnostic and nondiagnostic faces for each familiar face and at each level of feature intensity. To ensure that diagnostic and nondiagnostic faces produced the expected effect for each one of the four identities, we fitted a linear mixed effects model (i.e. fitlme, Matlab 2016b) to the data of each identity separately, using Wilkinson’s formulae:

$$\text{Performance} \sim 1 + \text{Face Type} + \text{Task Type} + \text{Amplification} \\ + (\text{Face Type} + \text{Task Type} + \text{Amplification} - 1 | \text{Subject})$$

The model had fixed factors of Face Type (i.e. diagnostic vs. nondiagnostic), Feature Amplification (i.e. 0.33, 0.67, 1, 1.33, 1.67) and Generalization Task (i.e. 3 views plus an age change and a sex change) as explanatory variables and participants’ response variability as random factor. From this model, we can infer whether or not the fixed factors generalized beyond the specific participant sample, separately for each identity.

We tested the specified fixed effect factor (i.e. using ANOVA, Matlab 2016b), using the Satherwith approximation to compute the approximate degrees of freedom. We found for each identity a higher identification performance with diagnostic than nondiagnostic faces (see Figure 4B), and the performance increased with amplification (an effect of Feature Amplification). The Generalization Task effect was significant for ‘Mary’ and ‘Stephany’ and not for ‘John’ and ‘Peter’. Supplementary Table 3 to 6 report the full statistics of our fixed effects, for each identity.

To further test the prediction effect of Face Type we built a null model that excludes this factor:

$$\text{Performance} \sim 1 + \text{Task Type} + \text{Amplification} + (\text{Task Type} + \text{Amplification} - 1 | \text{Subject})$$

For each identity, we compared the original and null model with a likelihood ratio (i.e. LR). Performance was significantly better explained by the original model (with Face Type) than the null model (without Face Type). For ‘Mary’, LR statistic = 603.72.135, $p < 0.001$; for ‘Stephany’, LR statistic = 39.516, $p < 0.001$; for ‘John’, LR

647 statistic = 205.67, $p < 0.001$; for 'Peter', LR statistic = 214.34, $p < 0.001$. See
648 Supplementary Table 3 to 6 for the full statistical analysis.

649 We also found a significant interaction effect between Face Type and
650 Amplification, by fitting a linear mixed effect model with this interaction included as an
651 effect factor (see Supplementary Methods, Linear Mixed Effect Model of Face Type
652 by Amplification Interaction, and Supplementary Table 7).

653 **Data Availability.** Data is available in Mendeley Data with identifier
654 <http://dx.doi.org/10.17632/nyt677xwfm.1> [50].

655 **Code Availability.** Analysis scripts are available in Mendeley Data with identifier
656 <http://dx.doi.org/10.17632/nyt677xwfm.1> [50].

657

658 REFERENCES

- 659 1. Bar, M. (2009). The proactive brain: memory for predictions. *Philos T R Soc B* 364,
660 1235-1243.
- 661 2. Bar, M., Kassam, K.S., Ghuman, A.S., Boshyan, J., Schmid, A.M., Dale, A.M.,
662 Hamalainen, M.S., Marinkovic, K., Schacter, D.L., Rosen, B.R., et al. (2006). Top-down
663 facilitation of visual recognition. *Proc Natl Acad Sci U S A* 103, 449-454.
- 664 3. Ullman, S., Assif, L., Fetaya, E., and Harari, D. (2016). Atoms of recognition in human
665 and computer vision. *P Natl Acad Sci USA* 113, 2744-2749.
- 666 4. Harel, A., Kravitz, D.J., and Baker, C.I. (2014). Task context impacts visual object
667 processing differentially across the cortex. *Proc Natl Acad Sci U S A* 111, E962-971.
- 668 5. O'Toole, A.J. (2011). Cognitive and Computational Approaches to Face Recognition In
669 The Oxford Handbook of Face Perception, G. Rhodes, A. Calder, M. Johnson and J.V.
670 Haxby, eds., pp. 15 -30.
- 671 6. Tsao, D.Y., and Livingstone, M.S. (2008). Mechanisms of face perception. *Annu Rev*
672 *Neurosci* 31, 411-437.
- 673 7. Rosch, E., and Mervis, C.B. (1975). Family Resemblances - Studies in Internal
674 Structure of Categories. *Cognitive Psychol* 7, 573-605.
- 675 8. Ahumada, A., and Lovell, J. (1971). Stimulus Features in Signal Detection. *J Acoust*
676 *Soc Am* 49, 1751-&.
- 677 9. Yu, H., Garrod, O.G.B., and Schyns, P.G. (2012). Perception-driven facial expression
678 synthesis. *Comput Graph-Uk* 36, 152-162.
- 679 10. Lee, D.D., and Seung, H.S. (1999). Learning the parts of objects by non-negative
680 matrix factorization. *Nature* 401, 788-791.
- 681 11. Lee, H., and Kuhl, B.A. (2016). Reconstructing Perceived and Retrieved Faces from
682 Activity Patterns in Lateral Parietal Cortex. *J Neurosci* 36, 6069-6082.
- 683 12. Nestor, A., Plaut, D.C., and Behrmann, M. (2016). Feature-based face
684 representations and image reconstruction from behavioral and neural data. *Proc*
685 *Natl Acad Sci U S A* 113, 416-421.

- 686 13. Chang, C.H., Nemrodov, D., Lee, A.C.H., and Nestor, A. (2017). Memory and
687 Perception-based Facial Image Reconstruction. *Sci Rep-Uk* 7.
- 688 14. Turk, M., and Pentland, A. (1991). Eigenfaces for recognition. *J Cogn Neurosci* 3, 71-
689 86.
- 690 15. Cootes, T.F., Edwards, G.J., and Taylor, C.J. (2001). Active appearance models. *Ieee T*
691 *Pattern Anal* 23, 681-685.
- 692 16. Blanz, V., and Vetter, T. (1999). A morphable model for the synthesis of 3D faces.
693 *Comp Graph*, 187-194.
- 694 17. Rhodes, G., and Jeffery, L. (2006). Adaptive norm-based coding of facial identity.
695 *Vision Res* 46, 2977-2987.
- 696 18. O'Toole, A.J., Castillo, C.D., Parde, C.J., Hill, M.Q., and Chellappa, R. (2018). Face
697 Space Representations in Deep Convolutional Neural Networks. *Trends Cogn Sci* 22,
698 794 - 809.
- 699 19. Young, A.W., and Burton, A.M. (2018). Are We Face Experts? *Trends in Cognitive*
700 *Sciences* 22, 100-110.
- 701 20. White, D., Phillips, P.J., Hahn, C.A., Hill, M., and O'Toole, A.J. (2015). Perceptual
702 expertise in forensic facial image comparison. *Proc Biol Sci* 282.
- 703 21. Eger, E., Schweinberger, S.R., Dolan, R.J., and Henson, R.N. (2005). Familiarity
704 enhances invariance of face representations in human ventral visual cortex: fMRI
705 evidence. *Neuroimage* 26, 1128-1139.
- 706 22. Jenkins, R., White, D., Van Montfort, X., and Burton, A.M. (2011). Variability in
707 photos of the same face. *Cognition* 121, 313-323.
- 708 23. Gosselin, F., and Schyns, P.G. (2002). RAP: a new framework for visual categorization.
709 *Trends Cogn Sci* 6, 70-77.
- 710 24. Schyns, P.G. (1998). Diagnostic recognition: task constraints, object information, and
711 their interactions. *Cognition* 67, 147-179.
- 712 25. Palmeri, T.J., Wong, A.C.N., and Gauthier, I. (2004). Computational approaches to
713 the development of perceptual expertise. *Trends in Cognitive Sciences* 8, 378-386.
- 714 26. Burton, A.M., Schweinberger, S.R., Jenkins, R., and Kaufmann, J.M. (2015).
715 Arguments Against a Configural Processing Account of Familiar Face Recognition.
716 *Perspect Psychol Sci* 10, 482-496.
- 717 27. Erens, R.G., Kappers, A.M., and Koenderink, J.J. (1993). Perception of local shape
718 from shading. *Percept Psychophys* 54, 145-156.
- 719 28. Phong, B.T. (1975). Illumination for Computer Generated Pictures. *Commun Acn* 18,
720 311-317.
- 721 29. Liu, Z.L. (1996). Viewpoint dependency in object representation and recognition.
722 *Spatial Vision* 9, 491-521.
- 723 30. Schyns, P.G., Goldstone, R.L., and Thibaut, J.P. (1998). The development of features
724 in object concepts. *Behav Brain Sci* 21, 1-17; discussion 17-54.
- 725 31. Mangini, M.C., and Biederman, I. (2004). Making the ineffable explicit: estimating
726 the information employed for face classifications. *Cognitive Sci* 28, 209-226.
- 727 32. Baxter, M.G. (2009). Involvement of medial temporal lobe structures in memory and
728 perception. *Neuron* 61, 667-677.

- 729 33. Xu, T., Zhan, J., Garrod, O.G.B., Torr, P.H.S., Zhu, S.C., Ince, R.A., and Schyns, P.G.
730 (2018). Deeper Interpretability of Deep Networks. ArXiv.
- 731 34. Leopold, D.A., O'Toole, A.J., Vetter, T., and Blanz, V. (2001). Prototype-referenced
732 shape encoding revealed by high-level aftereffects. *Nat Neurosci* 4, 89-94.
- 733 35. Leopold, D.A., Bondar, I.V., and Giese, M.A. (2006). Norm-based face encoding by
734 single neurons in the monkey inferotemporal cortex. *Nature* 442, 572-575.
- 735 36. Chang, L., and Tsao, D.Y. (2017). The Code for Facial Identity in the Primate Brain.
736 *Cell* 169, 1013-1028 e1014.
- 737 37. Zhan, J., Ince, R.A.A., van Rijsbergen, N., and Schyns, P.G. (2019). Dynamic
738 Construction of Reduced Representations in the Brain for Perceptual Decision
739 Behavior. *Curr Biol* 29, 319-326 e314.
- 740 38. Kay, K.N., Naselaris, T., Prenger, R.J., and Gallant, J.L. (2008). Identifying natural
741 images from human brain activity. *Nature* 452, 352-U357.
- 742 39. Smith, F.W., and Muckli, L. (2010). Nonstimulated early visual areas carry
743 information about surrounding context. *P Natl Acad Sci USA* 107, 20099-20103.
- 744 40. Peirce, J.W. (2015). Understanding mid-level representations in visual processing. *J*
745 *Vis* 15, 5.
- 746 41. Kubilius, J., Wagemans, J., and Op de Beeck, H.P. (2014). A conceptual framework of
747 computations in mid-level vision. *Front Comput Neurosci* 8, 158.
- 748 42. Friston, K.J., and Kiebel, S. (2009). Predictive coding under the free-energy principle.
749 *Philos T R Soc B* 364, 1211-1221.
- 750 43. Clark, A. (2013). Whatever next? Predictive brains, situated agents, and the future of
751 cognitive science. *Behavioral and Brain Sciences* 36, 181-204.
- 752 44. Gosselin, F., and Schyns, P.G. (2003). Superstitious perceptions reveal properties of
753 internal representations. *Psychol Sci* 14, 505-509.
- 754 45. Smith, M.L., Gosselin, F., and Schyns, P.G. (2012). Measuring Internal
755 Representations from Behavioral and Brain Data. *Current Biology* 22, 191-196.
- 756 46. Nestor, A., Plaut, D.C., and Behrmann, M. (2011). Unraveling the distributed neural
757 code of facial identity through spatiotemporal pattern analysis. *P Natl Acad Sci USA*
758 108, 9998-10003.
- 759 47. Gobbini, M.I., Gors, J.D., Halchenko, Y.O., Rogers, C., Guntupalli, J.S., Hughes, H., and
760 Cipolli, C. (2013). Prioritized Detection of Personally Familiar Faces. *PLoS One* 8,
761 e66620.
- 762 48. van Belle, G., Ramon, M., Lefevre, P., and Rossion, B. (2010). Fixation patterns during
763 recognition of personally familiar and unfamiliar faces. *Front Psychol* 1, 20.
- 764 49. Ramon, M., Vizioli, L., Liu-Shuang, J., and Rossion, B. (2015). Neural microgenesis of
765 personally familiar face recognition. *Proc Natl Acad Sci U S A* 112, E4835-4844.
- 766 50. Zhan, J., Garrod, O.G., Van Rijsbergen, N., and Schyns, P. Modelling Face Memory
767 Reveals Task-Generalizable Representations. *Mendeley Data*
768 <http://dx.doi.org/10.17632/nyt677xwfm.1> (2019)

769

770

771 **Acknowledgements.** P.G.S. received support from the Wellcome Trust (Senior
772 Investigator Award, UK; 107802) and the Multidisciplinary University Research
773 Initiative/Engineering and Physical Sciences Research Council (USA, UK; 172046-
774 01). The funders had no role in study design, data collection and analysis, decision to
775 publish or preparation of the manuscript.

776 **Competing interests.** The authors declare no competing interests.

777 **Author Contributions.** J.Z., N.VR and P.G.S. designed the research; O.G. and
778 P.G.S. developed the Generative Model of 3D Faces; J.Z. performed the research;
779 J.Z. and N.VR. analysed the data; and J.Z., N.VR. and P.G.S. wrote the paper.

Figure 1. Reverse correlating mental representations of familiar faces. (A) Task. Illustrative experimental trial with 6 randomly generated face identities. We instructed participants to use their memory to select the face most similar to a familiar identity (here, 'Mary') and then to rate the similarity of the selected face (purple frame) to their memory of 'Mary' (purple pointer). (B) Generative Model of 3D face identity (GMF). In its forward computation flow (see identity modelling solid arrow), the General Linear Model (GLM) decomposes a 3D, textured face (e.g. 'Jane' or 'Tom') into a non-identity face shape average capturing the categorical factors of face sex, ethnicity, age and their interactions plus a separate component that defines the identity of the face (illustrated by the 3D shape decomposition; 2D texture, not illustrated, is independently and similarly decomposed). Heat maps indicate the 3D shape deviations that define 'Jane' and 'Tom' in the GMF in relation to their categorical averages. In the reverse flow (see dashed arrow of identity generation), we can randomize the 3D shape identity component (and 2D texture component, not illustrated here), add the categorical average of 'Jane' (or 'Tom') and generate random faces, each with a unique identity that share all other categorical face information with 'Jane' and 'Tom.'

Figure 2. Contents of mental representations of familiar faces. (A) Mental representation of 'Mary' (a typical participant). *Ground truth:* 3D vertex positions deviate both Inward (-) and Outward (+) from the categorical average to objectively define the shape of each familiar face identity. Greyscale values reported on the flanking faces color-code the normalized magnitudes of inward and outward deviations from the categorical average. *Mental representation:* Inward and Outward colored faces highlight the individual 3D vertices whose position faithfully deviate from the categorical average in the GMF ($p < 0.05$, two-sided). Blue to red colors represent the normalized magnitudes of their deviations. *2D scatter plots:* Scatter plots indicate the relationship between each vertex deviation in the ground truth (normalized scale on the X-axis) and the corresponding vertex in the memory representation (normalized scale on the Y-axis). The white diagonal line provides the reference of veridical mental representation in the GMF—i.e. a hypothetical numerical correspondence between each shape vertex position in the ground truth face and in the mental representation of the same face. White dots indicate vertices that were not faithfully represented. (B) Mental Representations (group results). Same caption as Figure 2A, except that the colormap now reflects the number of participants ($N = 14$) who faithfully represented this particular shape vertex.

Figure 3. NNMF multivariate and compact representations. A. NNMF representations of faithful 3D vertices across the mental representations of participants. The x-axis heatmap presents each NNMF component, where colors indicate the relative weight of each shape vertex in the component (normalized by maximum weight across components). Boxplots on the y-axis show the loading of each NNMF component on the faithful representations ($N = 14$, one per participant) of each familiar identity ($N = 4$ familiar identities), with colored boxes indicating above 0.1 threshold loading for NNMF components. In boxplots, the bottom (vs. top) edges indicate the 25th (vs. 75th) percentile of the distribution; the whiskers cover the +2.7

822 standard deviation; the larger central circle indicates the median; the outliers are plotted in
823 smaller circle outside the whiskers. B. Diagnostic and nondiagnostic components for each
824 familiar identity. Heat maps in the left column show the diagnostic component for each
825 familiar identity; heat maps in the right column show the complementary nondiagnostic
826 components.

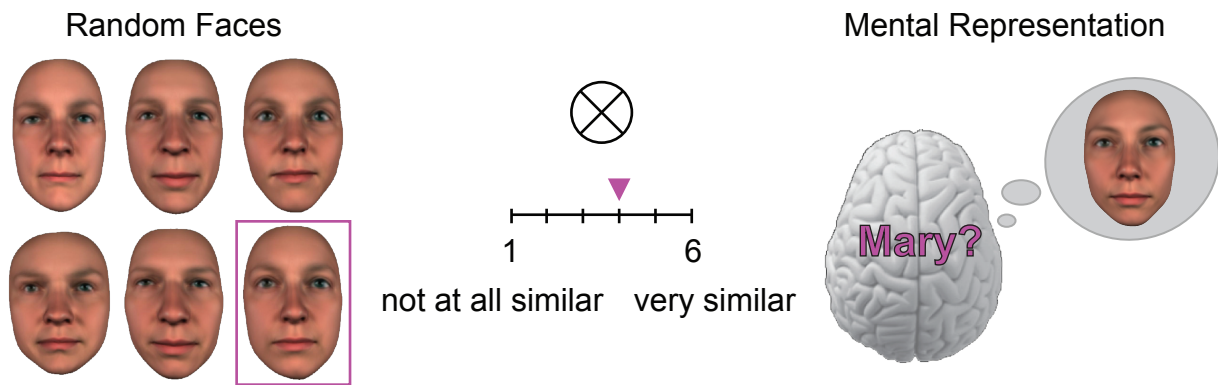
827

828 **Figure 4. Generalization of performance across tasks.** (A) Diagnostic and nondiagnostic
829 Faces. *Left panel:* The red background map shows the multivariate diagnostic components of
830 faithful 3D shape representation of 'Mary'; the grey background map shows the nondiagnostic
831 complement (1 - diagnostic components). *Middle panel:* Faces synthesized with increasing
832 amplification (0.33 to 1.67) of the diagnostic (top) vs. nondiagnostic (bottom) components.
833 *Right panel:* For each synthesized face, we changed its viewpoint (30° left and 30° right), age
834 (80 years old) and sex, shown here for faces synthesized at amplification = 1. (B) Task
835 Performance. For each condition of generalization (row) and familiar identity (column), 2D
836 plots show the median identification performance computed across 12 validators (y-axes) for
837 faces synthesized with the diagnostic (red curves) and nondiagnostic (grey curves) faces, at
838 different levels of amplification of the multivariate components (x-axes). Shadowed regions
839 indicate median absolute deviations (MAD) of identification performance. Abbreviations: Diag
840 = Diagnostic, Nondiag = Nondiagnostic.

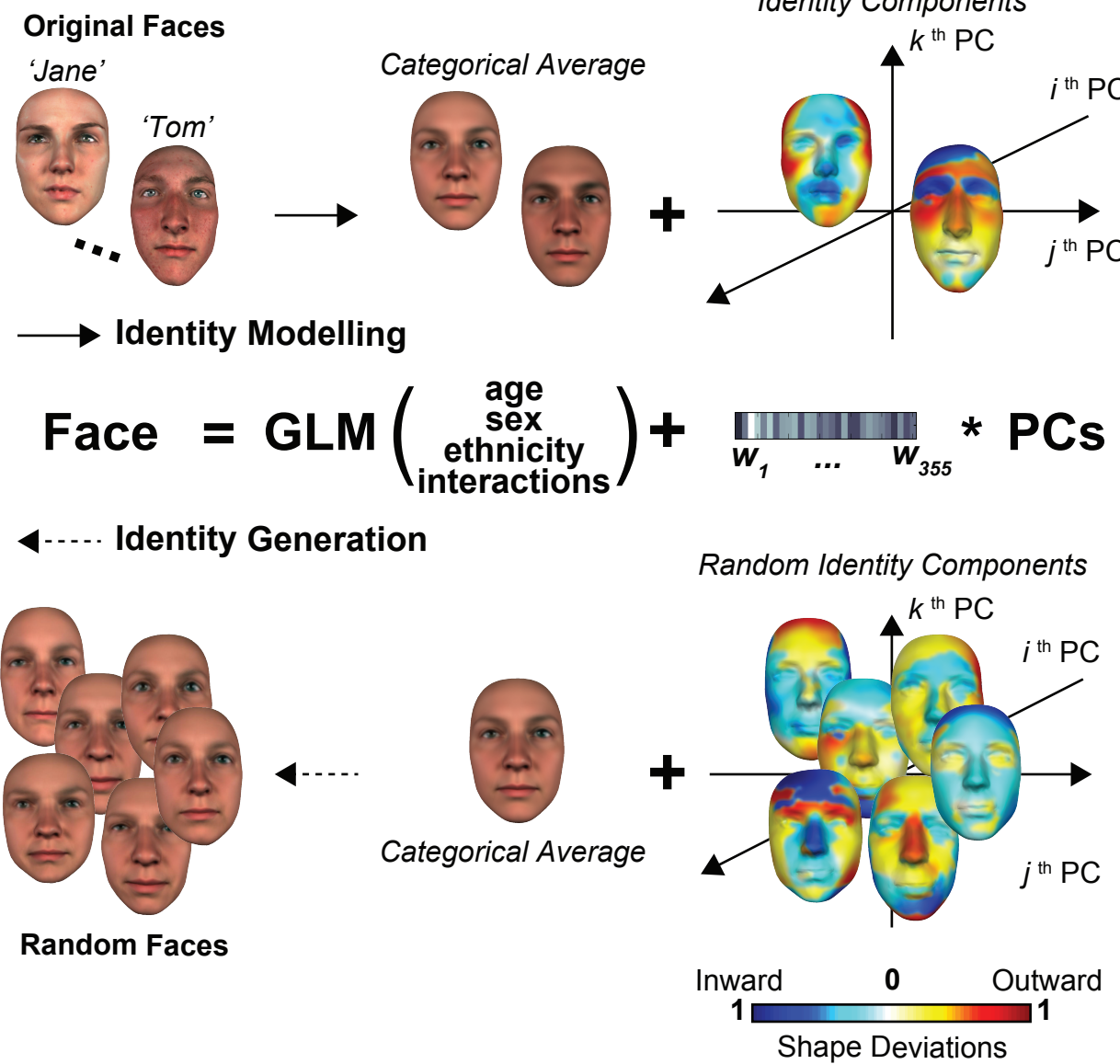
841

842

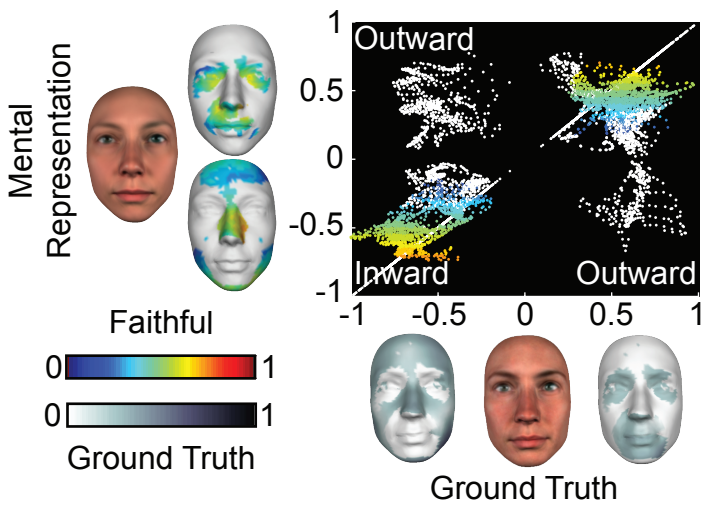
A. Reverse Correlation Task



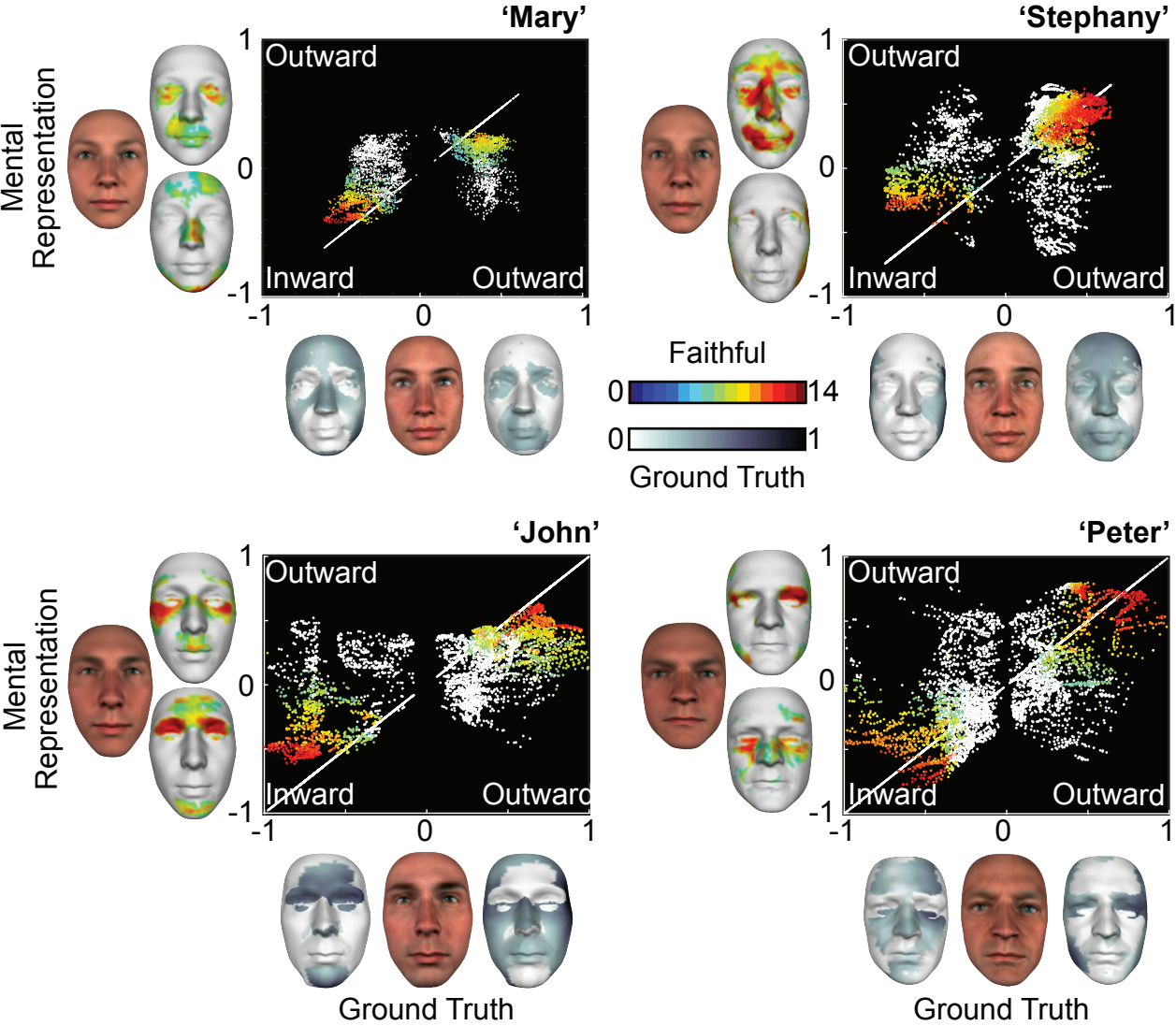
B. Generative Model of 3D Face Identity



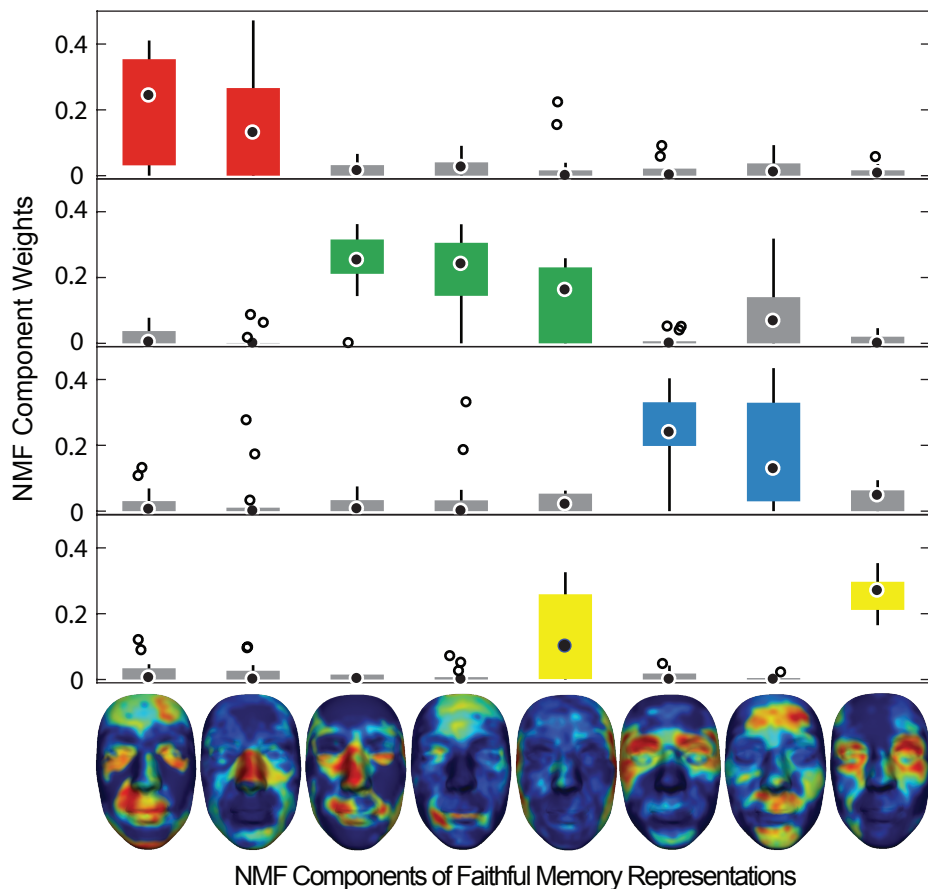
A. Mental Representation of ‘Mary’ (One Participant)



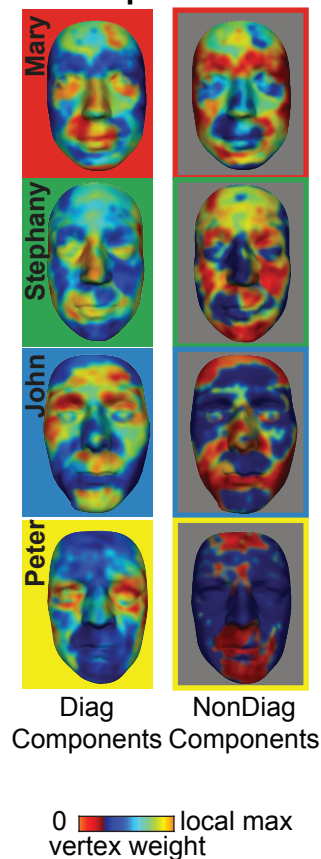
B. Mental Representations (Group Results)



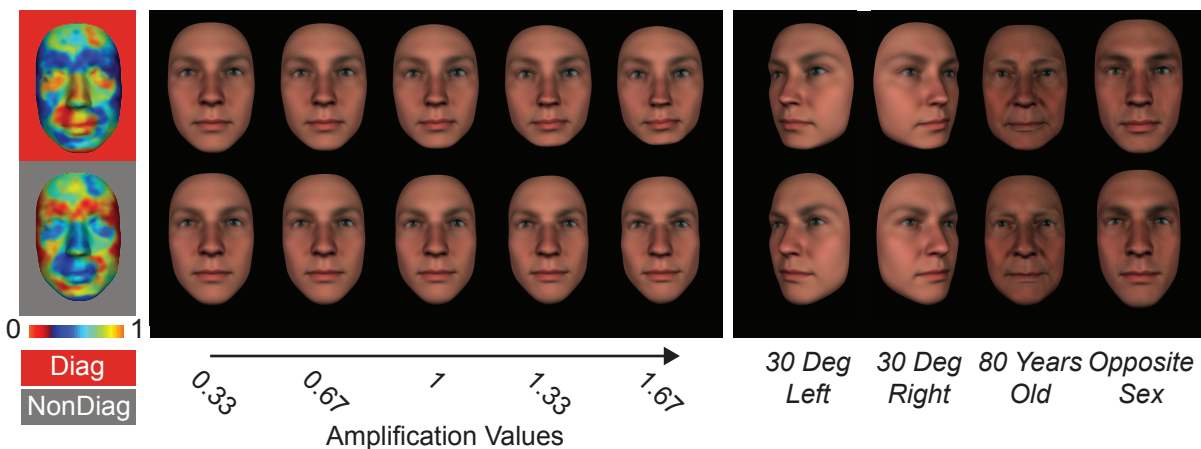
A. Multivariate NMF Representations of Faithful 3D Vertices



B. Diag vs. NonDiag Components



A. Diagnostic and Nondiagnostic Faces



B. Identification Performance of Diagnostic and Nondiagnostic Faces

